

Data Analysis with PANDAS

CHEAT SHEET

CREATED BY: ARIANNE COLTON AND SEAN CHEN

DATA STRUCTURES

SERIES (1D)

One-dimensional array-like object containing an array of data (of any **Numpy** data type) and an associated array of data labels, called its **"index"**. If index of data is not specified, then a default one consisting of the integers 0 through N-1 is created.

Create Series	<code>series1 = pd.Series([1, 2], index = ['a', 'b'])</code> <code>series1 = pd.Series(dict1) *</code>
Get Series Values	<code>series1.values</code>
Get Values by Index	<code>series1['a']</code> <code>series1[['b', 'a']]</code>
Get Series Index	<code>series1.index</code>
Get Name Attribute (None is default)	<code>series1.name</code> <code>series1.index.name</code>
** Common Index Values are Added	<code>series1 + series2</code>
Unique But Unsorted	<code>series2 = series1.unique()</code>

- * Can think of Series as a fixed-length, ordered dict. Series can be substituted into many functions that expect a dict.
- ** Auto-align differently-indexed data in arithmetic operations

DATAFRAME (2D)

Tabular data structure with ordered collections of columns, each of which can be different value type. Data Frame (DF) can be thought of as a dict of Series.

Create DF (from a dict of equal-length lists or Numpy arrays)	<code>dict1 = {'state': ['Ohio', 'CA'], 'year': [2000, 2010]}</code> <code>df1 = pd.DataFrame(dict1)</code> # columns are placed in sorted order <code>df1 = pd.DataFrame(dict1, index = ['row1', 'row2'])</code> # specifying index <code>df1 = pd.DataFrame(dict1, columns = ['year', 'state'])</code> # columns are placed in your given order
* Create DF (from nested dict of dicts) The inner keys as row indices	<code>dict1 = {'col1': {'row1': 1, 'row2': 2}, 'col2': {'row1': 3, 'row2': 4}}</code> <code>df1 = pd.DataFrame(dict1)</code>

Get Columns and Row Names	<code>df1.columns</code> <code>df1.index</code>
Get Name Attribute (None is default)	<code>df1.columns.name</code> <code>df1.index.name</code>
Get Values	<code>df1.values</code> # returns the data as a 2D ndarray, the dtype will be chosen to accommodate all of the columns
** Get Column as Series	<code>df1['state']</code> or <code>df1.state</code>
** Get Row as Series	<code>df1.ix['row2']</code> or <code>df1.ix[1]</code>
Assign a column that doesn't exist will create a new column	<code>df1['eastern'] = df1.state == 'Ohio'</code>
Delete a column	<code>del df1['eastern']</code>
Switch Columns and Rows	<code>df1.T</code>

- * Dicts of Series are treated the same as Nested dict of dicts.
- ** Data returned is a 'view' on the underlying data, NOT a copy. Thus, any in-place modifications to the data will be reflected in df1.

PANEL DATA (3D)

Create Panel Data : (Each item in the Panel is a DF)

<code>import pandas_datareader.data as web</code> <code>panel1 = pd.Panel({'stk': web.get_data_yahoo(stk, '1/1/2000', '1/1/2010') for stk in ['AAPL', 'IBM']})</code> # panel1 Dimensions : 2 (item) * 861 (major) * 6 (minor)	
"Stacked" DF form : (Useful way to represent panel data)	<code>panel1 = panel1.swapaxes('item', 'minor')</code> <code>panel1.ix[:, '6/1/2003', :].to_frame() *</code> => Stacked DF (with hierarchical indexing **): # Open High Low Close Volume Adj-Close # major minor # 2003-06-01 AAPL # IBM # 2003-06-02 AAPL # IBM

DATA STRUCTURES CONTINUED

- * DF has a `"to_panel()"` method which is the inverse of `"to_frame()"`.
- ** Hierarchical indexing makes N-dimensional arrays unnecessary in a lot of cases. Aka prefer to use Stacked DF, not Panel data.

INDEX OBJECTS

Immutable objects that hold the axis labels and other metadata (i.e. axis name)

- * i.e. Index, MultiIndex, DatetimeIndex, PeriodIndex
- * Any sequence of labels used when constructing Series or DF internally converted to an Index.
- * Can functions as fixed-size set in addition to being array-like.

HIERARCHICAL INDEXING

Multiple index levels on an axis : A way to work with higher dimensional data in a lower dimensional form.

MultiIndex :	<code>series1 = Series(np.random.randn(6), index = [['a', 'a', 'a', 'b', 'b', 'b'], [1, 2, 3, 1, 2, 3]])</code> <code>series1.index.names = ['key1', 'key2']</code>
Series Partial Indexing	<code>series1['b']</code> # Outer Level <code>series1[:, 2]</code> # Inner Level
DF Partial Indexing	<code>df1['outerCol3', 'InnerCol2']</code> Or <code>df1['outerCol3']['InnerCol2']</code>

Swapping and Sorting Levels

Swap Level (level interchanged) *	<code>swapSeries1 = series1.swaplevel('key1', 'key2')</code>
Sort Level	<code>series1.sortlevel(1)</code> # sorts according to first inner level

Common Ops : Swap and Sort **	<code>series1.swaplevel(0, 1).sortlevel(0)</code> # the order of rows also change
-------------------------------	--

- * The order of the rows do not change. Only the two levels got swapped.
- ** Data selection performance is much better if the index is sorted starting with the outermost level, as a result of calling `sortlevel(0)` or `sort_index()`.

Summary Statistics by Level

Most stats functions in DF or Series have a **"level"** option that you can specify the level you want on an axis.

Sum rows (that have same 'key2' value)	<code>df1.sum(level = 'key2')</code>
Sum columns ..	<code>df1.sum(level = 'col3', axis = 1)</code>

- * Under the hood, the functionality provided here utilizes panda's **"groupby"**.

DataFrame's Columns as Indexes

DF's **"set_index"** will create a new DF using one or more of its columns as the index.

New DF using columns as index	<code>df2 = df1.set_index(['col3', 'col4']) * ‡</code> # col3 becomes the outermost index, col4 becomes inner index. Values of col3, col4 become the index values.
-------------------------------	---

- * **"reset_index"** does the opposite of **"set_index"**, the hierarchical index are moved into columns.
- ‡ By default, 'col3' and 'col4' will be removed from the DF, though you can leave them by option: `'drop = False'`.

MISSING DATA

Python	NaN = np.nan (not a number)
Pandas *	NaN or python built-in None mean missing/NA values

* Use `pd.isnull()`, `pd.notnull()` or `series1/df1.isnull()` to detect missing data.

FILTERING OUT MISSING DATA

`dropna()` returns with ONLY non-null data, source data NOT modified.

<code>df1.dropna()</code>	# drop any row containing missing value
<code>df1.dropna(axis = 1)</code>	# drop any column containing missing values

<code>df1.dropna(how = 'all')</code>	# drop row that are all missing
<code>df1.dropna(thresh = 3)</code>	# drop any row containing < 3 number of observations

FILLING IN MISSING DATA

<code>df2 = df1.fillna(0)</code>	# fill all missing data with 0
<code>df1.fillna(inplace = True)</code>	# modify in-place
Use a different fill value for each column :	
<code>df1.fillna({'col1': 0, 'col2': -1})</code>	
Only forward fill the 2 missing values in front :	
<code>df1.fillna(method = 'ffill', limit = 2)</code>	
i.e. for column1, if row 3-6 are missing. so 3 and 4 get filled with the value from 2, NOT 5 and 6.	